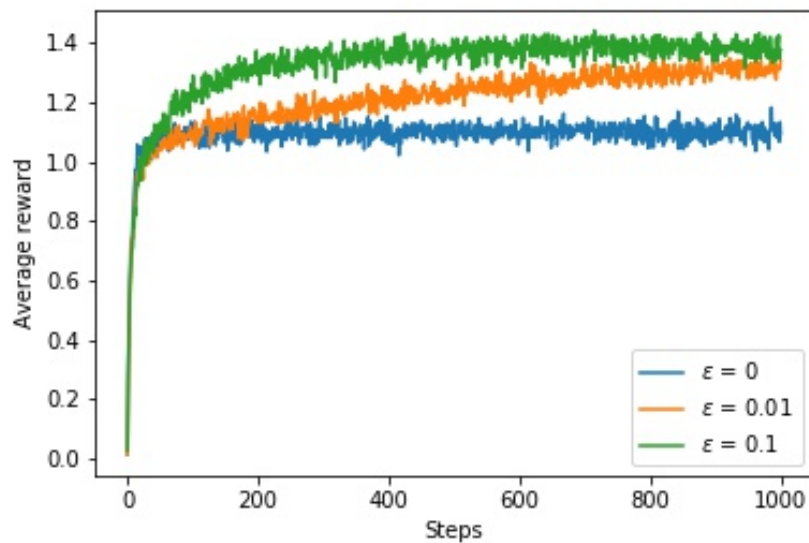


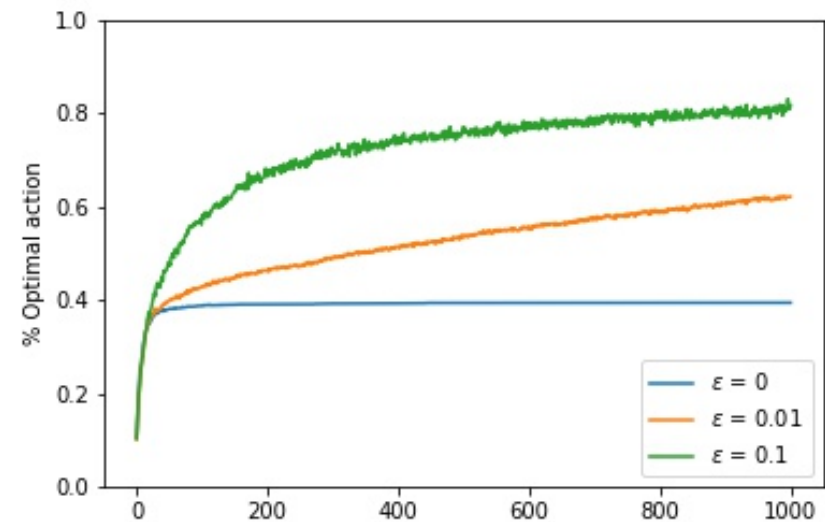
RL Assignment 1

Srivatsava Kesanupalli MT18054

Exercise 2.5 - The 10-armed bandit problem for $\epsilon = 0, 0.01, 0.1$



Average reward



Optimal action

Figure 2.3 - $Q_1 = 5$ with $\epsilon = 0$ vs. $Q_1 = 0$ with $\epsilon = 0.1$

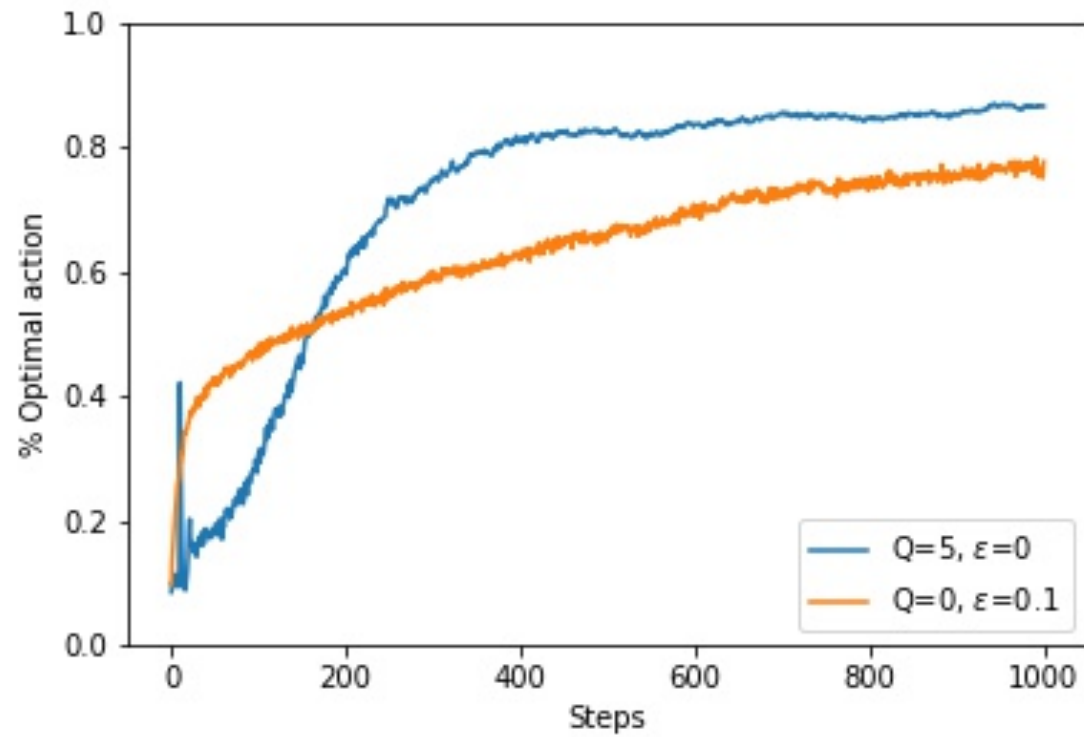
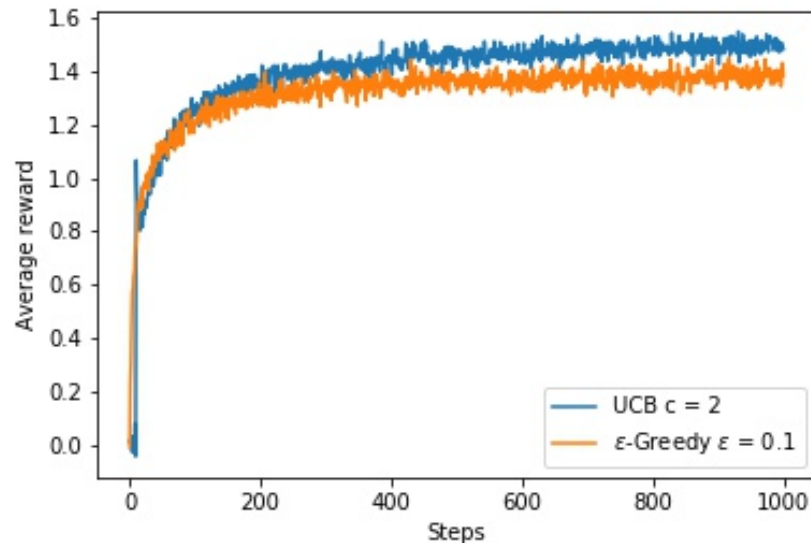
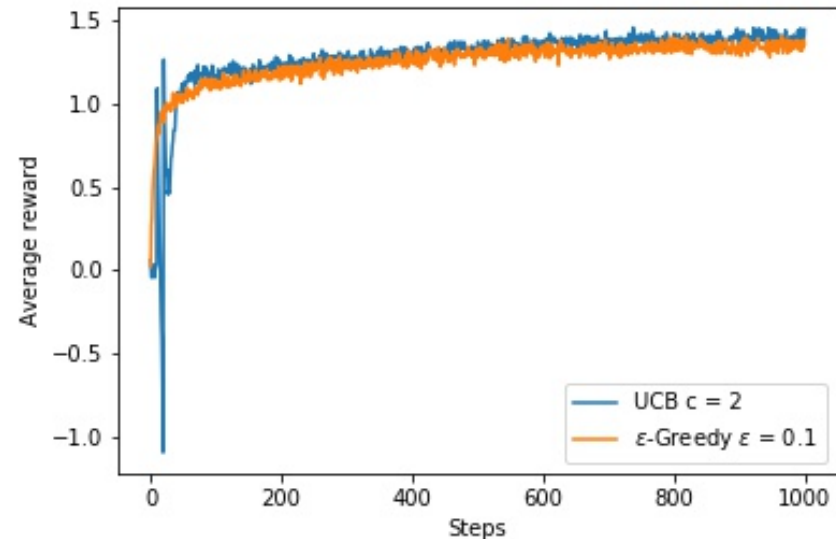


Fig 2.3

UCB vs. ϵ -Greedy



UCB stationary vs ϵ -Greedy



UCB non-stationary vs ϵ -Greedy

UCB performs better than ϵ -Greedy as the error term in the estimate penalises for not exploring an action. The performance can be seen for both stationary as well as non-stationary cases.

